

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2000-035965

(43)Date of publication of application : 02.02.2000

(51)Int.Cl.

G06F 17/30
G06T 7/00
// G06T 1/00

(21)Application number : 10-203583

(71)Applicant : NIPPON TELEGR & TELEPH CORP
<NTT>

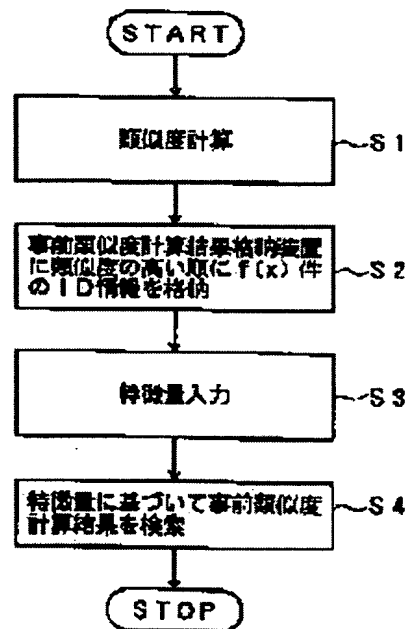
(22)Date of filing : 17.07.1998

(72)Inventor : AKAMA HIROKI
SATO MICHİYOSHI
MITSUI KAZUYOSHI
KUSHIMA KAZUHIKO(54) METHOD AND DEVICE FOR RETRIEVING SIMILAR FEATURE QUANTITY AND STORAGE
MEDIUM STORING RETRIEVAL PROGRAM OF SIMILAR FEATURE QUANTITY

(57)Abstract:

PROBLEM TO BE SOLVED: To perform fast retrieval by retrieving a pre-similarity calculation result storage device based on retrieval key feature quantity when feature quantity in a database is given as the retrieval key feature quantity and returning a pre-similarity calculation result as a retrieval result.

SOLUTION: The whole feature quantities are preliminarily made keys, similar calculation in a database is performed and other feature quantities and similarity are calculated (S1). ID information for upper rank $f(x)$ matters is stored in a pre-similarity calculation result storage device in order of similarity with a similarity sequence or together with similarity value as occasion demands (S2). And, when feature quantity in the database is given as retrieval key feature quantity (S3), the pre-similarity calculation result storage device is retrieved based on the retrieval key feature quantity and a pre-similarity calculation result is returned as a retrieval result (S4). According to this method, it is possible to very fast perform retrieval even if the number of the entire database feature quantity data pieces is large.



LEGAL STATUS

[Date of request for examination] 21.11.2000

[Date of sending the examiner's decision of rejection] 25.11.2003

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

Copyright (C); 1998,2003 Japan Patent Office

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開2000-35965

(P 2 0 0 0 - 3 5 9 6 5 A)

(43) 公開日 平成12年2月2日 (2000.2.2)

| (51) Int. Cl. ⁷ | 識別記号 | F I | テーマコード (参考) |
|-----------------------------|------|----------------|-------------|
| G06F 17/30 | | G06F 15/40 370 | G 5B050 |
| G06T 7/00 | | 15/401 310 | D 5B075 |
| // G06T 1/00 | | 15/403 350 | C 5L096 |
| | | 15/70 460 | B |
| | | 15/62 330 | G |
| 審査請求 未請求 請求項の数12 O L (全17頁) | | | |

(21) 出願番号 特願平10-203583

(22) 出願日 平成10年7月17日 (1998.7.17)

(71) 出願人 000004226

日本電信電話株式会社

東京都千代田区大手町二丁目3番1号

(72) 発明者 赤間 浩樹

東京都新宿区西新宿三丁目19番2号 日本
電信電話株式会社内

(72) 発明者 佐藤 路恵

東京都新宿区西新宿三丁目19番2号 日本
電信電話株式会社内

(74) 代理人 100070150

弁理士 伊東 忠彦

最終頁に続く

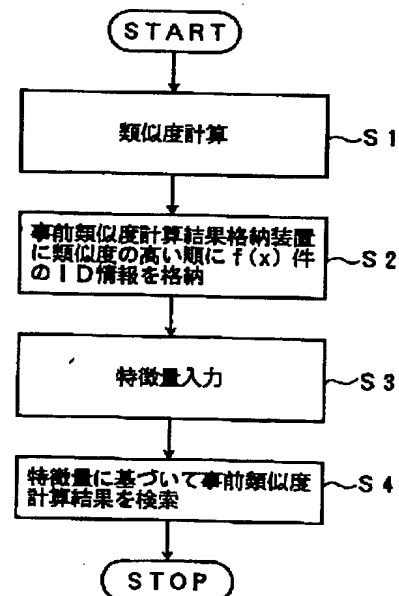
(54) 【発明の名称】 類似特徴量の検索方法及び装置及び類似特徴量の検索プログラムを格納した記憶媒体

(57) 【要約】

【課題】 事前類似度計算結果がディスク上または、部分的にディスク上に存在するような複雑な構造を持っていても、さらに、データベース全体の特徴量データ件数が多くとも高速な検索が可能な類似特徴量の検索方法及び装置及び類似特徴量の検索プログラムを格納した記憶媒体を提供する。

【解決手段】 本発明は、予めデータベース内の全ての特徴量をキーとし、データベース内における類似計算を行い、他の特徴量との類似度を計算し、類似度の高い順に上位 $f(x)$ 件分の ID 情報を、類似度順付で、あるいは、必要に応じて該類似度の値と共に、事前類似度計算結果格納装置に格納しておき、検索キー特徴量としてデータベース内の特徴量が与えられた場合、該検索キー特徴量に対する事前類似度計算結果を検索結果として返却する。

本発明の原理を説明するための図



【特許請求の範囲】

【請求項 1】 マルチメディアデータに対する類似検索システムやテキストの類似検索システムに用いられる類似特徴量の検索方法において、
予めデータベース内の全ての特徴量をキーとし、
前記データベース内における類似計算を行い、他の特徴量との類似度を計算し、
前記類似度の高い順に上位 $f(x)$ 件分の ID 情報を、
類似度順付で、あるいは、必要に応じて該類似度の値と共に、事前類似度計算結果格納装置に格納しておき、
検索キー特徴量として前記データベース内の特徴量が与えられた場合、該検索キー特徴量に基づいて前記事前類似度計算結果格納装置を検索して、事前類似度計算結果を検索結果として返却することを特徴とする類似特徴量の検索方法。

【請求項 2】 前記検索キー特徴量として前記データベース内に存在することが分からない特徴量が与えられた場合に、
前記特徴量に最も類似する前記データベース内の特徴量を最近傍検索装置により検索し、
検索結果の特徴量に基づいて前記事前類似度計算結果格納装置を検索して、事前類似度計算結果を検索結果として返却する請求項 1 記載の類似特徴量の検索方法。

【請求項 3】 特徴量データの追加がある場合に、追加されたデータに関しては、追加特徴量データ管理装置で管理を行い、
検索キー特徴量が与えられた場合には、
前記事前類似度計算結果格納装置を検索した結果と、前記追加特徴量データ管理装置からの検索結果を類似度順にマージした結果を検索結果として返却する請求項 1 乃至 2 記載の類似特徴量の検索方法。

【請求項 4】 特徴量データの追加がある場合に、前記追加特徴量データ管理装置内の特徴量データ数が特定値 t を越えた後に、または、特定間隔の時間経過を含むタイミングにより、追加特徴量データを含めたデータベース内の全てのデータに関し、前記事前類似度計算結果の再計算を検索を行う処理とは、独立にまたは、並列に行い、
計算が完了した時点で、前記事前類似度計算結果、及び前記追加特徴量データ管理装置のデータの更新を行う請求項 3 記載の類似特徴量の検索方法。

【請求項 5】 マルチメディアデータに対する類似検索システムやテキストの類似検索システムに用いられる類似特徴量の検索装置であって、
全ての特徴量をキーとするデータベースと、
前記データベース内における類似計算を行い、他の特徴量との類似度を計算する類似度計算手段と、
前記類似度計算手段で求められた前記類似度の高い順に上位 $f(x)$ 件分の ID 情報を、類似度順付で、あるいは、必要に応じて該類似度の値と共に格納する、事前類

似度計算結果格納手段と、

検索キー特徴量として前記データベース内の特徴量が与えられた場合、該検索キー特徴量に対する事前類似度計算結果を前記事前類似度計算結果格納手段を検索することにより取得して、検索結果として返却する事前類似度計算結果検索手段とを有することを特徴とする類似特徴量の検索装置。

【請求項 6】 前記検索キー特徴量として前記データベース内に存在することが分からない特徴量が与えられた場合に、前記特徴量に最も類似する前記データベース内の特徴量を検索する最近傍検索手段を更に有し、
前記事前類似度計算結果検索手段は、前記最近傍検索手段の検索結果の特徴量に対する事前類似度計算結果を検索結果として返却する請求項 5 記載の類似特徴量の検索装置。

【請求項 7】 特徴量データの追加がある場合に、追加されたデータに関して管理する追加特徴量データ管理手段と、
検索キー特徴量が与えられた場合には、前記事前類似度計算結果格納手段からの結果と、前記追加特徴量データ管理手段からの検索結果を類似度順にマージした結果を検索結果として返却するマージ手段を有する請求項 5 乃至 6 記載の類似特徴量の検索装置。

【請求項 8】 特徴量データの追加がある場合に、前記追加特徴量データ管理手段内の特徴量データ数が特定値 t を越えた後に、または、特定間隔の時間経過を含むタイミングにより、追加特徴量データを含めたデータベース内の全てのデータに関し、前記事前類似度計算結果検索手段とは独立または、並列に事前類似度計算を行う再計算手段と、

前記再計算手段の計算が完了した時点で、前記事前類似度計算結果格納手段、及び前記追加特徴量データ管理手段のデータの更新を行う更新手段を有する請求項 7 記載の類似特徴量の検索装置。

【請求項 9】 マルチメディアデータに対する類似検索システムやテキストの類似検索システムに用いられる類似特徴量の検索プログラムを格納した記憶媒体であって、
データベース内の全ての特徴量をキーとするデータベース内における類似計算を行い、他の特徴量との類似度を計算する類似度計算プロセスと、
前記類似度計算プロセスで求められた前記類似度の高い順に上位 $f(x)$ 件分の ID 情報を、類似度順付で、あるいは、必要に応じて該類似度の値と共に事前類似度計算結果格納手段に格納する事前類似度計算結果格納制御プロセスと、検索キー特徴量として前記データベース内の特徴量が与えられた場合、該検索キー特徴量に対する事前類似度計算結果を前記事前類似度計算結果格納手段を検索することにより取得して、検索結果として返却する事前類似度計算結果検索プロセスとを有することを特

徴とする類似特徴量の検索プログラムを格納した記憶媒体。

【請求項 10】 前記検索キー特徴量として前記データベース内に存在することが分からない特徴量が与えられた場合に、前記特徴量に最も類似する前記データベース内の特徴量を検索する最近傍検索プロセスを更に有し、前記事前類似度計算結果検索プロセスは、前記最近傍検索プロセスの検索結果の特徴量に対する事前類似度計算結果を検索結果として返却する請求項 9 記載の類似特徴量の検索プログラムを格納した記憶媒体。

【請求項 11】 検索キー特徴量が与えられた場合には、前記事前類似度計算結果格納手段からの結果と、特徴量データの追加がある場合に、追加されたデータに関して管理する追加特徴量データ管理手段からの検索結果を類似度順にマージした結果を検索結果として返却するマージプロセスを有する請求項 8 乃至 10 記載の類似特徴量の検索プログラムを格納した記憶媒体。

【請求項 12】 特徴量データの追加がある場合に、前記追加特徴量データ管理手段内の特徴量データ数が特定値 t を越えた後に、または、特定間隔の時間経過を含むタイミングにより、追加特徴量データを含めたデータベース内の全てのデータに関し、前記事前類似度計算結果検索プロセスとは独立または、並列に事前類似度計算を行う再計算プロセスと、前記再計算プロセスの計算が完了した時点で、前記事前類似度計算結果格納手段、及び前記追加特徴量データ管理手段のデータの更新を行う更新プロセスを含む請求項 11 記載の類似特徴量の検索プログラムを格納した記憶媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】 本発明は、類似特徴量の検索方法及び装置及び類似特徴量の検索プログラムを格納した記憶媒体に係り、特に、画像、映像、モーション、音楽、音声などのマルチメディアデータに対する類似検索システムの実現やテキストの類似検索システムに用いられる類似特徴量の検索方法及び装置及び類似特徴量の検索プログラムを格納した記憶媒体に関する。詳しくは、インターネット上の画像のように、大量で、その量が日々増加するような対象に対し、高速な類似検索を実現することに用いるための類似特徴量の検索方法及び装置及び類似特徴量の検索プログラムを格納した記憶媒体に関する。

【0002】

【従来の技術】 最初に多次元特徴量データについて説明する。画像検索、音楽検索などに代表される検索は、従来の RDBMS が対象としてきた一致検索や範囲検索とは異なり、多次元特徴量（次元数は 1 以上）の類似検索である。

【0003】 ここで、一致検索とは、データベース内の

列に対する検索キー値が与えられた時、それと一致する値を持つ全行、または、行 ID を検索結果とする検索をいう。範囲検索とは、データベース内の列に対し、検索キーとしての値と共に、検索条件として大小関係が与えられ、その関係を満足する値を持つデータベース内の全行、または、行 ID を検索結果とする検索をいう。

【0004】 一方、類似検索とは、1 次元以上の多次元特徴量（一般に単に特徴量と呼ぶ。複数の数値からなるためベクトルと呼ぶこともある）をデータベース格納の対象とし、検索キーとして与えられた特徴量キーに対し、その特徴量間の距離等を計算することにより類似度を求め、最も類似度の高い順に上位 $f(x)$ 件の行を求めるような検索を行う。

【0005】 特徴量としては、画像や音楽等マルチメディア情報の内容特徴などのこともあるし、地図座標のこともあるし、テキスト内のキーワードの重みのこともある。類似検索は、範囲検索の対象を 1 次元の値から多次元のベクトル値に拡張した場合に似ているが、その返却値の考え方が異なり、範囲検索の場合は、範囲条件が明確に指定され、その条件を満たす行は全て検索結果となるものの、類似検索の場合は、一般には、類似の高い順に上位 $f(x)$ 件を返すという指定が用いられる。

【0006】 以下の明細書中の説明において、上位 $f(x)$ 件と記述した場合に、それは抽象化された値を示しており、単に、利用者が指定した特定の値 k 、システムが予め持つ特定の値 k 、また、最大 k や最小 k 、データベース内の全データ数、利用者、システムまたは、データベースの状態から得られる情報等から計算によって求められた利用者または、システムまたは、データベースの状態から得られる情報等から計算によって求められた値のように、別の手段で計算される閾値 k でもよい。また、図等で 1 つのフローチャート内に複数の $f(x)$ という表記があっても、それらは独立な値を持つてもよい。

【0007】 図 10 は、類似検索の例を説明するための図である。この例の特徴量は 2 次元で、データベース内には 6 件の特徴量データが登録されている。この利用者から与えられた検索キー特徴量 (0.5, 0.6) を入力した場合、データベース内の各特徴量とのユークリッド距離を計算し、その距離の近い順に並べ替え、その中の上位何件かを検索結果として利用者に返却する。

【0008】 次に、高速化について説明する。最も単純な類似検索では、検索キー特徴量とデータベース内の全特徴量データとの類似度計算が検索実行時に行われる。ところで、この特徴量が 1 次元の場合には、従来の R データベース MS の範囲検索で利用されていたような手法 (B+Tree 等) を用いることで高速検索が可能になる。

【0009】 しかし、類似検索では、一般には特徴量は 2 以上の次元数となるため、上記の手法は利用できない

い。そこで、その高速化のための索引手法には以下のような手法が用いられる。図11、図12は、R-treeの例を示しており、図11は、従来のR-treeの特徴量空間分割を説明するための図であり、図12は、従来のT-treeの木構造を説明するための図である。構成される木の各ノードは、どの次元で分割したかという情報と、その範囲の情報を持つ。各分割は、その中に含まれる特徴量点の個数が同程度になるように調整されている。図13、図14は、PR-quadtreesの例を示しており、図13は、従来のPR-quadtreesの特徴量空間分割を説明するための図であり、図14は、従来のPR-quadtreesの木構造を説明するための図である。空間は常にX-Y平面で4つに分割され、分割後の領域に指定個数以上の特徴量点が存在する場合は、さらに4分割が行われていく。

【0010】それぞれに関し、各種の改良バージョンが提案されているが（参考：Volker Gaedo, Multidimensional Access Methods）、一般には、大量のデータに対しても、その木を平衡状態を維持するR-tree、及びその改良系が高速性、汎用性に優れている。本発明では、これらの多次元空間を分割し、木状にした構造をもつ索引を木状索引と呼ぶことにする。

【0011】図15は、従来の木状索引を使った類似検索のフローチャートであり、木状索引の構築時の流れ、及び木状索引を使った類似検索時の流れを示している。索引構築時は、特徴量の部分集合をデータベース全体の特徴量とし（ステップ10）、特徴量数または、リンク数が（木のノード内数）以上であれば、特徴量の部分集合に対して以下の処理を行う（ステップ11）、分割基準を決定し（ステップ12）、分割基準によって特徴量集合をn個の分割し（ステップ13）、個々の集合に対し、再帰的に繰り返す（ステップ14）。再帰終了の場合には階層的な分類結果を索引として登録する（ステップ15）。

【0012】検索実行時は、検索キー特徴量を入力し（ステップ20）、与えられた特徴量がどの分類に相当するか、分類基準に従って索引を辿る（ステップ21）。

【0013】

【発明が解決しようとする課題】しかしながら、上記従来の木状索引による高速な検索手法には以下のような問題がある。最初に高次元数特徴量データでの検索速度における観点から説明する。従来の木状検索手法は、特徴量データ数の増加に対して、その検索速度の増加を抑えることを主な目的としている。つまり、特徴量データを木構造で管理することで、特徴量データ間の比較階数をlogのオーダーとし、特徴量データ数の増加に対する検索性能を維持する。

【0014】しかし、これら従来の木状索引構造は、次元数の増加に対しては考慮されておらず、例えば、R-

tree等では、数次元程度で最も威力を発揮し、20次元を越えるとその性能はかなり悪くなることが知られている。これは、地理情報等、低次元の応用には充分だが、マルチメディア情報等の高次元の応用には不十分である。

【0015】次に、高度類似基準への対処における観点から説明する。従来の索引手法は、マンハッタン距離（市街地距離）やユークリッド距離のように数学的に距離の公理を満たす単純な類似度基準を想定している。これらの類似度基準により、事前にデータベース内のデータ間の関係を各次元軸をもとにクラスタリングした場合には、そのクラスタリング結果空間の中で近いデータ同士は、その元となる類似度基準でも近いという性質があり、事前に木状索引の作成が可能になる。しかし、その類似度基準が与えられた検索キーデータに依存し、各次元を元に事前にクラスタリングすることが意味をなさない場合、例えば、ヒストグラム、インターセクション（参考：Maichael J. Swain, Indexing Via Color Histogram）や、非対象類似度（参考：赤間、オブジェクトの類似度算出方法及び類似オブジェクト検索装置、特願平9-060999）といった、マルチメディア情報の特徴量に合った複雑な類似度基準には対応できない。なお、本明細書では、距離をより一般化した用語として類似度を用いている。

【0016】次に、近傍順検索時の検査速度の観点から説明する。木状に構成された索引では、最近傍データを検出するのは容易である。しかし、一般的な類似検索においては、最も類似するものだけを検索するに留まらず、似ている順に上位f(x)件の結果を求めることが多い。その場合、木状に管理されたデータにおいては、木の枝や葉を順に辿り、候補の中のデータに関して、再度、類似度の計算を行う必要がある。また、これは、特徴量データ数が増加し、データがメモリ上ではなく、ディスク上にある場合は、かなりの速度低下要因となる。

【0017】最後に、実装法の観点から説明する。木の平衡状態を維持する等、アルゴリズムが複雑で実装が困難である。本発明は、上記の点に鑑みなされたもので、事前類似度計算結果がディスク上または、部分的にディスク上に存在するような複雑な構造を持っていても、さらに、データベース全体の特徴量データ件数が多くとも高速な検索が可能な類似特徴量の検索方法及び装置及び類似特徴量の検索プログラムを格納した記憶媒体を提供することを目的とする。

【0018】

【課題を解決するための手段】図1は、本発明の原理を説明するための図である。本発明（請求項1）は、マルチメディアデータに対する類似検索システムやテキストの類似検索システムに用いられる類似特徴量の検索方法において、予めデータベース内の全ての特徴量をキーとし、データベース内における類似計算を行い、他の特徴

量との類似度を計算し(ステップ1)、類似度の高い順に上位 $f(x)$ 件分のID情報を、類似度順付で、あるいは、必要に応じて該類似度の値と共に、事前類似度計算結果格納装置に格納しておき(ステップ2)、検索キー特徴量としてデータベース内の特徴量が与えられた場合(ステップ3)、該検索キー特徴量に基づいて事前類似度計算結果格納装置を検索して、事前類似度計算結果を検索結果として返却する(ステップ4)。

【0019】本発明(請求項2)は、検索キー特徴量としてデータベース内に存在することが分からない特徴量が与えられた場合に、特徴量に最も類似するデータベース内の特徴量を最近傍検索装置により検索し、検索結果の特徴量に基づいて事前類似度計算結果格納装置を検索して、事前類似度計算結果を検索結果として返却する。

【0020】本発明(請求項3)は、特徴量データの追加がある場合に、追加されたデータに関しては、追加特徴量データ管理装置で管理を行い、検索キー特徴量が与えられた場合には、事前類似度計算結果格納装置を検索した結果と、追加特徴量データ管理装置からの検索結果を類似度順にマージした結果を検索結果として返却する。

【0021】本発明(請求項4)は、特徴量データの追加がある場合に、追加特徴量データ管理装置内の特徴量データ数が特定値 t を越えた後に、または、特定間隔の時間経過を含むタイミングにより、追加特徴量データを含めたデータベース内の全てのデータに関し、事前類似度計算結果の再計算を検索を行う処理とは、独立にまたは、並列に行い、計算が完了した時点で、事前類似度計算結果、及び追加特徴量データ管理装置のデータの更新を行う。

【0022】図2は、本発明の原理構成図である。本発明(請求項5)は、マルチメディアデータに対する類似検索システムやテキストの類似検索システムに用いられる類似特徴量の検索装置であって、全ての特徴量をキーとするデータベース10と、データベース10における類似計算を行い、他の特徴量との類似度を計算する類似度計算手段20と、類似度計算手段20で求められた類似度の高い順に上位 $f(x)$ 件分のID情報を、類似度順付で、あるいは、必要に応じて該類似度の値と共に格納する、事前類似度計算結果格納手段30と、検索キー特徴量としてデータベース10内の特徴量が与えられた場合、該検索キー特徴量に対する事前類似度計算結果を事前類似度計算結果格納手段30を検索することにより取得して、検索結果として返却する事前類似度計算結果検索手段40とを有する。

【0023】本発明(請求項6)は、検索キー特徴量としてデータベース10内に存在することが分からない特徴量が与えられた場合に、特徴量に最も類似するデータベース10内の特徴量を検索する最近傍検索手段を更に有し、事前類似度計算結果検索手段40は、最近傍検索

手段の検索結果の特徴量に対する事前類似度計算結果を検索結果として返却する。

【0024】本発明(請求項7)は、特徴量データの追加がある場合に、追加されたデータに関して管理する追加特徴量データ管理手段と、検索キー特徴量が与えられた場合には、事前類似度計算結果格納手段30からの結果と、追加特徴量データ管理手段からの検索結果を類似度順にマージした結果を検索結果として返却するマージ手段を有する。

【0025】本発明(請求項8)は、特徴量データの追加がある場合に、追加特徴量データ管理手段内の特徴量データ数が特定値 t を越えた後に、または、特定間隔の時間経過を含むタイミングにより、追加特徴量データを含めたデータベース10内の全てのデータに関し、事前類似度計算結果検索手段40とは独立または、並列に事前類似度計算を行う再計算手段と、再計算手段の計算が完了した時点で、事前類似度計算結果格納手段30、及び追加特徴量データ管理手段のデータの更新を行う更新手段を有する。

【0026】本発明(請求項9)は、マルチメディアデータに対する類似検索システムやテキストの類似検索システムに用いられる類似特徴量の検索プログラムを格納した記憶媒体であって、データベース内の全ての特徴量をキーとするデータベース内における類似計算を行い、他の特徴量との類似度を計算する類似度計算プロセスと、類似度計算プロセスで求められた類似度の高い順に上位 $f(x)$ 件分のID情報を、類似度順付で、あるいは、必要に応じて該類似度の値と共に事前類似度計算結果格納手段に格納する事前類似度計算結果格納制御プロセスと、検索キー特徴量としてデータベース内の特徴量が与えられた場合、該検索キー特徴量に対する事前類似度計算結果を事前類似度計算結果格納手段を検索することにより取得して、検索結果として返却する事前類似度計算結果検索プロセスとを有する。

【0027】本発明(請求項10)は、検索キー特徴量としてデータベース内に存在することが分からない特徴量が与えられた場合に、特徴量に最も類似するデータベース内の特徴量を検索する最近傍検索プロセスを更に有し、事前類似度計算結果検索プロセスは、最近傍検索プロセスの検索結果の特徴量に対する事前類似度計算結果を検索結果として返却する。

【0028】本発明(請求項11)は、検索キー特徴量が与えられた場合には、事前類似度計算結果格納手段からの結果と、特徴量データの追加がある場合に、追加されたデータに関して管理する追加特徴量データ管理手段からの検索結果を類似度順にマージした結果を検索結果として返却するマージプロセスを有する。本発明(請求項12)は、特徴量データの追加がある場合に、追加特徴量データ管理手段内の特徴量データ数が特定値 t を越えた後に、または、特定間隔の時間経過を含むタイミン

グにより、追加特徴量データを含めたデータベース内の全てのデータに関し、事前類似度計算結果検索プロセスとは独立または、並列に事前類似度計算を行う再計算プロセスと、再計算プロセスの計算が完了した時点で、事前類似度計算結果格納手段、及び追加特徴量データ管理手段のデータの更新を行う更新プロセスを含む。

【0029】上記のように、本発明は、類似度の高い順に上位 $f(x)$ 件分のID情報に類似度順が付与された事前類似度計算結果に対する最近傍検索処理は、既にデータベース内に存在する値に対する一致検索となるため、その索引方法としては、従来のB-Tree、B+Tree、ハッシュ等のごく一般的な（容易な）手法を利用することで実現できる。また、近傍順検索については、事前に計算してある結果をそのまま返却するだけの処理となるため、その結果が例えば、ディスク上にあろうが、部分的にディスク上に存在するような複雑な構造を持っていようが、非常に高速に検索が可能になる。

【0030】また、検索実行時に次元数に依存する類似度計算を行うことがないため、次元数の増加に対しても性能が劣化することが少なく、高速である。さらに、索引の構造の中に距離に依存した部分がないので、特殊な類似性基準にも対応できる。また、データベース内の特徴量のみを対象とする場合には、最も類似する特徴量は、必ず自分自身であるため、一般には出力するか否かについてシステムに依存するが、データベース外特徴量を対象とする場合には、通常、最も類似する特徴量を出力する必要がある。

【0031】また、特徴データの追加がある場合でも事前類似度計算結果格納手段からの結果と追加特徴量データ管理手段に格納されている検索結果を類似度順にマージして、上位 $f(x)$ 件を検索結果として出力することができる。これにより、追加データのあるシステムの場合においても事前類似度計算結果を索引として利用することが可能となる。

【0032】

【発明の実施の形態】以下の説明において、特徴量データをデータベース内に存在する特徴量（これをデータベース内特徴量と呼ぶ）と、データベース内に存在しない特徴量（これをデータベース外特徴量と呼ぶ）の2種類に分けて考える。例えば、類似画像検索システムにおいて、データベース外特徴量を検索キーとして利用する例としては、スケッチ入力画像を検索キーとする場合や、デジタルカメラ画像を検索キーとする場合がある。

【0033】一方、データベース内部特徴量のIDを検索キーとして利用する例としては、システムが利用者にランダムに提示した画像を検索キーとする場合や、キーワード検索等の他の手法で検索した画像を検索キーとする場合や、一度検索した結果を利用してナビゲーション的に繰り返し検索する場合などがある。類似検索の索引の処理を、検索キー特徴量に最も類似する特徴量を求め

る処理である最近傍検索と、最近傍検索で求めた特徴量から順に近い特徴量を求めていく近傍順検索の2つの処理を分けて考えると、データベース外特徴量を検索キーとする類似検索では、最近傍検索と近傍順検索の両方が必要であり、データベース内特徴量のIDを検索キーとする類似検索では、近傍順検索のみ必要である。

【0034】なお、検索キーとしてデータベース内特徴量そのものが与えられた場合においても、一致検索によってデータベース内特徴量IDに変換することが可能であるため、最近傍検索は必要ない。本発明では、主に近傍順検索の処理部分の高速化を対象とする。図3は、本発明の類似特徴量検索装置の構成を示す。

【0035】同図に示す類似特徴量検索装置は、データベース10、類似度計算部20、事前類似度計算結果格納部30、検索部40、検索キー入力部50、特徴量種別判定部55、出力部60、最近傍検索部70、追加特徴量データ管理部80、マージ部90から構成される。データベース10は、全ての特徴量をキーとして、ID情報及びデータと共に格納する。

【0036】類似度計算部20は、データベース10内における類似計算を行い、他の特徴量との類似度を計算し、類似度の高い順に上位 $f(x)$ 件分のID情報に類似度順を付与してデータベース10に事前類似度計算結果格納部30に格納する。必要によっては、当該類似度の値と共に、事前類似度計算結果格納部30に格納する。

【0037】事前類似度計算結果格納部30は、類似度計算部20により求められた類似度計算結果（類似度順、類似度が付与されたID情報）を格納する。検索部40は、検索キー入力部50により与えられた検索キー特徴量として特徴量が与えられると、事前類似度計算結果格納部30を検索して、上位 $f(x)$ 件を検索結果として出力部60より出力する。

【0038】検索キー入力部50は、検索キー特徴量として特徴量を入力する。特徴量種別判定部55は、検索キー入力部50から入力された特徴量がデータベース10にあるか、データベース10外にあるかを判定する。出力部60は、検索部40、最近傍検索部70及びマージ部80で求められた検索結果を出力する。

【0039】最近傍検索部70は、検索キー特徴量として検索キー入力部50からデータベース10内に存在するか否かが分からない特徴量が与えられた場合には、それに最も類似するデータベース内の特徴量をR-tree等を用いて検索し、その結果の特徴量に対する事前類似度計算結果格納部30から検索して、上位 $f(x)$ 件を検索結果として返却する。

【0040】追加特徴量データ管理部80は、検索キー入力部50から入力された特徴量データを格納する。マージ部90は、検索部40が事前類似度計算結果格納部30から検索した検索結果と、追加特徴量データ管理部

80から検索した検索結果とをマージする。次に、上記の構成における動作を説明する。

【0041】図4は、本発明の検索構築時及び検索実行時の動作を示すフローチャートである。まず、最初に検索構築時の動作について説明する。

ステップ101) データベース10内における全特徴量に対して以下の処理を繰り返す。

【0042】ステップ102) 類似度計算部20は、データベース10内における類似度計算を行い、他の特徴量との類似度の計算を行い、類似度の高い順に上位 $f(x)$ 件分のID情報を、類似度順、類似度の値を求め

る。
ステップ103) 類似度計算部20により求められた結果を、特徴量または、そのIDをキーとして事前類似度計算結果格納部30に格納する。

【0043】次に、検索実行時の動作について説明する。

ステップ201) データベース10内の特徴量を検索キー特徴量として検索キー入力部50より入力される。

ステップ202) 検索部40は、入力された特徴量または、そのIDをキーとして確定検索方式により事前類似度計算結果格納部30に対して検索を行い、検索結果を取得する。

【0044】ステップ203) 検索部40は、上位 $f(x)$ 件分の結果を出力部60に出力する。このときの事前類似度計算結果に対する最近傍検索の処理は、すでにデータベース10内に存在する値に対する一致検索となるため、その検索方法としては、既存のB-Tree、B+Tree、ハッシュ等の一般的な(容易な)手法を利用することで実現できる。

【0045】また、近傍順検索については、事前に計算してある結果をそのまま返却するのみの処理となるため、その結果が例えば、データベース10上にあるが、部分的にデータベース10上に存在するような複雑な構造を持っていようが、非常に高速に検索が可能となる。また、検索実行時に次元数に依存する類似度計算を行うことがないため、次元数の増加に対しても性能が劣化することが少なく、高速である。

【0046】さらに、索引の構造の中に距離に依存した部分がないので、特殊な類似性基準にも対応できる。次に、最近傍検索の処理について説明する。データベース10内特徴量のみを対象とする場合には、最も類似する特徴量は必ずデータベース10内にあるため、一般には出力するか否かについてシステムに依存するが、データベース外特徴量を対象とする場合には、通常最も類似する特徴量を出力する必要がある。

【0047】図5は、本発明の最近傍検索の処理を示すフローチャートである。

ステップ301) 検索キー入力部50から検索キー特徴量として、データベース10内に存在することが分か

らない特徴量が与えられる。

ステップ302) 特徴量種別判定部55において、入力された特徴量がデータベース10にあるか、データベース10外にあるかを判定し、データベース10内にある場合にはステップ303に移行し、データベース10外にある場合にはステップ304に移行する。

【0048】ステップ303) 入力された特徴量がデータベース10外にある場合には、最近傍検索部70において、与えられた特徴がどの分類に相当するかを分類基準にしたがって検索を辿り、最も近い特徴量のIDを取得し、ステップ304に移行する。

ステップ304) 入力された特徴量がデータベース10内にある場合には、検索部40は、入力された特徴量または、最近傍検索部70により求められた特徴量のIDをキーとして確定検索方式により事前類似度計算結果格納部30に対して検索を行い、結果を出力部60に出力し、ステップ305に移行する。

【0049】ステップ305) 出力部60において、上位 $f(x)$ 件分の結果を出力する。次に、特徴量データの追加がある場合に対処する処理を説明する。図6は、本発明の特徴量データの追加がある場合の処理を示すフローチャートである。

【0050】ステップ401) まず、特徴量データの追加時の処理として、特徴量のデータの追加がある場合には、特徴量データの追加と索引の再構成を行い、追加特徴量データ管理部80に格納する。

ステップ501) 検索実行時の処理として、検索キー特徴量が検索キー入力部50から入力される。

【0051】ステップ502) 検索部40は、検索キー特徴量を用いて、事前類似計算結果格納部30から上位 $f(x)$ 件の類似検索を行う。

ステップ503) さらに、検索部40は、追加特徴量データ管理部80から上位 $f(x)$ 件以内の類似検索を行う。

ステップ504) マージ部90は、ステップ502とステップ503で求められた検索結果を距離順に整列させる。

【0052】ステップ505) 出力部60から上位 $f(x)$ 件分の結果を返却する。また、特徴量データの追加がある場合において、追加特徴量データ管理部80内の特徴量データ数が特定値 t を越えた後に、または、特定間隔の時間経過等のタイミングにより、その追加特徴量データを含むデータベース10内のデータに関し、検索部40における事前類似度計算結果格納部30による検索処理とは独立して、計算が完了した時点で、事前類似度計算結果格納部30及び追加特徴量データ管理部80のデータの更新を行う。

【0053】これにより、追加データのあるシステムの場合でも、事前類似度計算結果を索引として利用することが可能となる。但し、高度な類似度を用いる場合や高

次元特徴量を扱う場合や、データ数が少ない場合には、R - t r e e 等の木状索引を利用せずに、全件処理を行う方が望ましい。

【0054】なお、本明細書では、特徴量種が1種類の場合を想定して記述しているが、2種以上の特徴量が存在し、それらを独立に検索するような場合にも、複数の事前類似度計算結果を持つことで同様に適用できる。特徴量種別としては、画像の場合、色相、彩度、輝度、テクスチャ、大きさ等、画像オブジェクトの場合には、さらに、形、位置、傾き等多種存在する。

【0055】また、本明細書では、1種類の特徴量に対し、1種類の類似度基準を前提として記述しているが、複数の類似度基準（または、距離基準）を切り替えて検索を可能にするシステムに対応するため、事前類似度計算結果を類似度基準の種類数分だけ用意すればよい。

【0056】

【実施例】以下、図面と共に本発明の実施例を説明する。

【第1の実施例】図7は、本発明の第1の実施例の事前類似度計算結果の例を説明するための図であり、図8は、本発明の第1の実施例の総当たりによる事前類似度計算の例を説明するための図である。

【0057】図7、図8を用いて事前類似度計算結果の作成方法、及びその検索時の利用方法について説明する。まず、事前類似度計算結果を作成するため、データベース10内の全特徴量に対して以下の処理を繰り返す。始めに、キーをID1の(0.3, 0.3)とし、ID2～ID6までのデータを対象とした類似検索を行う。その結果が図8に示されている。この例では、最も簡単な実装の場合を想定し、「f(x)件」として全件（この場合6件）だったとした場合で、類似検索方法は、全ての組み合わせで類似度を計算した場合とし、類似度の高い順に並んだ6つのIDの結果、

ID1, ID2, ID4, ID5, ID3, ID6
を求め、図7のID1の事前類似度計算結果として登録している。

【0058】同様に、キーID2～キーID6までの処理を行った場合を図7に示す。なお、類似検索方法は、R - T r e e のような他の既存の類似検索用索引手法を用いた方法であっても構わない。また、事前類似度計算結果中には、必要に応じて、類似度等の情報を持つこともある。例えば、第3の実施例において後述するように、再度、類似度計算が必要な場合には、事前類似度計算結果として類似度情報を持つと効率がよい。

【0059】次に、事前類似度計算結果を使った検索の例を示す。検索キーとして与えられた特徴量がデータベース10内特徴量と分かる場合、通常、その情報はIDとして与えられ、IDを使って、事前類似度計算結果から、IDに割り当てられている事前類似度計算結果を得ることができる。しかし、もし、この段階で特徴量しか

与えられなかった場合でも、特徴量に対し、普通のB - t r e e 索引等が付与してあれば、単なる一致検索として、高速にそのIDを求めることができる。

【0060】なお、事前類似度計算結果に登録してあるIDの件数が、検索として要求され、検索された件数より少ない場合には、本発明では、上位f(x)件までの部分にしか機能せず、f(x)件の部分については、従来手法による類似順検索が必要になる。しかし、通常は、データベース作成時にアプリケーションとして利用する最大件数が決定できるため、それを越える個数のIDを事前類似度計算結果に用意しておけば問題ない。

【0061】【第2の実施例】本実施例では、与えられる検索キーが内部データベース特徴量と判断できない場合の例を示す。与えられた検索キーがデータベース内部にある特徴量と判断できない場合には、その特徴量データによる最近傍検索だけをR - T r e e のような他の従来手法を利用し、その後の近傍順検索については、本発明を利用する。

【0062】これは、例えば、図10に示す類似検索のように、検索キーとして(0.5, 0.6)が与えられた場合、その最近傍特徴量の(0.5, 0.5)を求めるまでは、従来手法を用い、その後、(0.5, 0.5)の近傍順検索では、そのIDに登録されている事前類似度計算結果のID4, ID3, ID1, ID5, ID2, ID6を検索結果とする。

【0063】なお、厳密な類似度順の結果を得たい場合には、再度、類似度計算を行い、整列をし直すものとする。

【第3の実施例】本実施例では、特徴量データに追加がある場合の処理を図7及び図9を用いて説明する。

【0064】図9は、本発明の第3の実施例の追加特徴量の管理とその類似検索の例を説明するための図である。本実施例において、事前類似度計算結果は図7に示すものとし、後に追加されたデータは、図9のように追加特徴量データ管理部80に格納される。この追加特徴量データ管理部80には、一般には従来の木状索引等が付与され、高速化される。また、データベース外特徴量を扱う場合の最近傍検索用索引と統合されることもある。

【0065】検索キー特徴量が与えられた場合には、事前類似度計算結果から上位f(x)件の類似検索結果を得、同時に、追加特徴量データ管理部80からも最大で上位f(x)件の類似検索結果を得る。このとき、それらの結果に類似度情報も付与しておき、その類似度で上位f(x)件の類似度データを作成し、それを類似検索結果とする。

【0066】なお、この検索キー特徴量がデータベースない特徴量の場合には、事前類似度計算結果から得られた上位f(x)件に対し、事前に計算された類似度を利

用することができるが、データベース外特徴量の場合には、類似度に関し、再計算が必要となる。例えば、検索キー特徴量が (0.5, 0.6) の場合、第 2 の実施例で示したように、事前類似度計算結果を使った検索結果は、

ID4, ID3, ID1, ID5, ID2, ID6

となり、その距離の再計算を行うと、

ID4, ID3, ID1, ID5, ID2, ID6

になる。

【0067】また、図 9 の追加特徴量データ管理部 80 から検索した結果は、

ID1, ID2, ID3

となり、これらを類似度順にマージすると、

内 ID4, 内 ID3, 追 ID1, 内 ID1, 内 ID6,

内 ID5, 追 ID2, 追 ID3, 内 ID2

となる。

【0068】但し、事前類似度計算結果内の ID は、

「内 ID」と、追加特徴量データ管理部 80 内の ID は、「追 ID」と記載し、区別した。よって、この内の上位 $f(x)$ 件を検索結果とすればよい。また、本発明は、上記の実施例に限定されることなく、図 3 に示す構成要件をプログラムとして構築し、類似特徴量検索装置として利用されるコンピュータに接続されるディスク装置や、フロッピーディスク、CD-ROM 等の可搬記憶媒体に格納しておき、本発明を実施する際に、インストールすることにより容易に本発明を実現できる。

【0069】なお、本発明は、上記の実施例に限定されることなく、特許請求の範囲内で種々変更・応用が可能である。

【0070】

【発明の効果】上述のように、本発明によれば、事前類似度計算結果情報が、たとえ、ディスク上にあろうが、部分的にディスク上に存在するような複雑な構造を持っていようが、データベース全体の特徴量データ件数が多からうが、非常に高速に検索ができる。

【0071】また、検索実行時に次元数に依存する類似度計算を行うことがないため、次元数の増加に対しても性能が劣化することが少なく、高速である。さらに、索引の構造の中に距離に依存した部分がないので、特殊な類似性基準にも対応できる。また、本発明は、近傍検索と組み合わせた、高速な類似検索が可能となる。

【0072】さらに、本発明は、追加の特徴量データが存在する場合にも、システムの構成が可能となる。

【図面の簡単な説明】

【図 1】本発明の原理を説明するための図である。

【図 2】本発明の原理構成図である。

【図 3】本発明の類似特徴量検索装置の構成図である。

【図 4】本発明の索引構築時及び検索実行時の動作を示すフローチャートである。

【図 5】本発明の最近傍検索の処理を示すフローチャートである。

【図 6】本発明の特徴量データの追加がある場合の処理を示すフローチャートである。

【図 7】本発明の第 1 の実施例の事前類似度計算結果の例を説明するための図である。

【図 8】本発明の第 1 の実施例の総当たりによる事前類似度計算の例を説明するための図である。

【図 9】本発明の第 3 の実施例の追加特徴量の管理とその類似検索の例を説明するための図である。

【図 10】類似検索を説明するための図である。

【図 11】従来の R-tree の特徴量空間分割を説明するための図である。

【図 12】従来の R-tree の木構造を説明するための図である。

【図 13】従来の PR-quadtrees の特徴量空間分割を説明するための図である。

【図 14】従来の PR-quadtrees の木構造を説明するための図である。

【図 15】従来の木状索引を使った類似検索のフローチャートである。

【符号の説明】

10 データベース

20 類似度計算手段、類似度計算部

30 事前類似度計算結果格納手段、事前類似度計算結果格納部

40 事前類似度計算結果検索手段、検索部

50 検索キー入力部

55 特徴量種別判定部

60 出力部

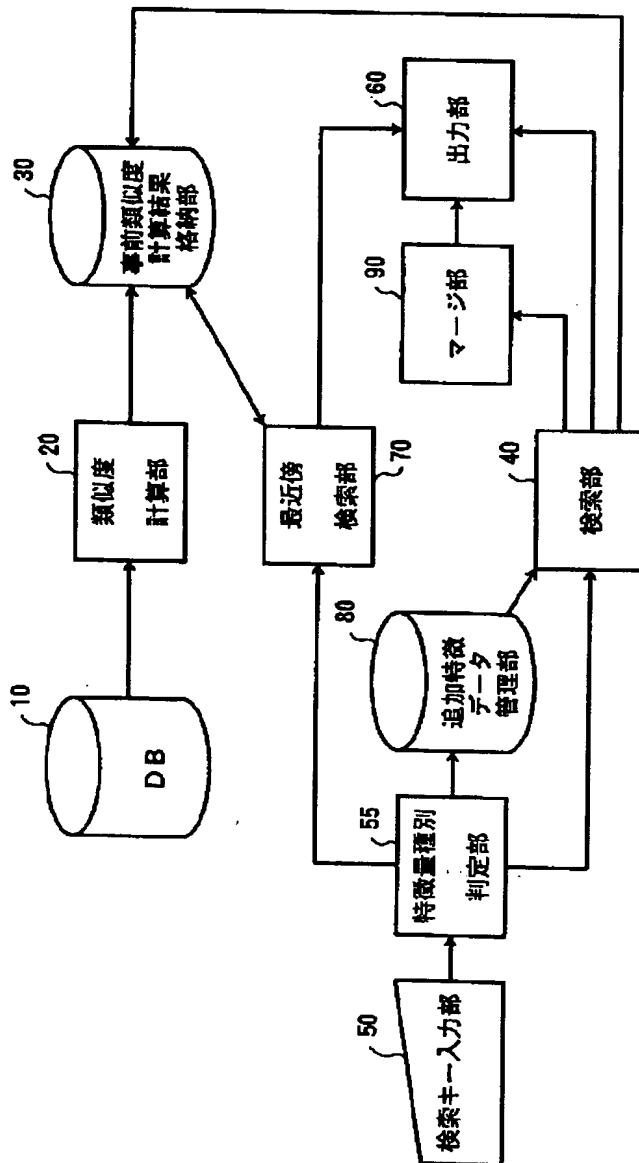
70 最近傍検索部

80 追加特徴量データ管理部

90 マージ部

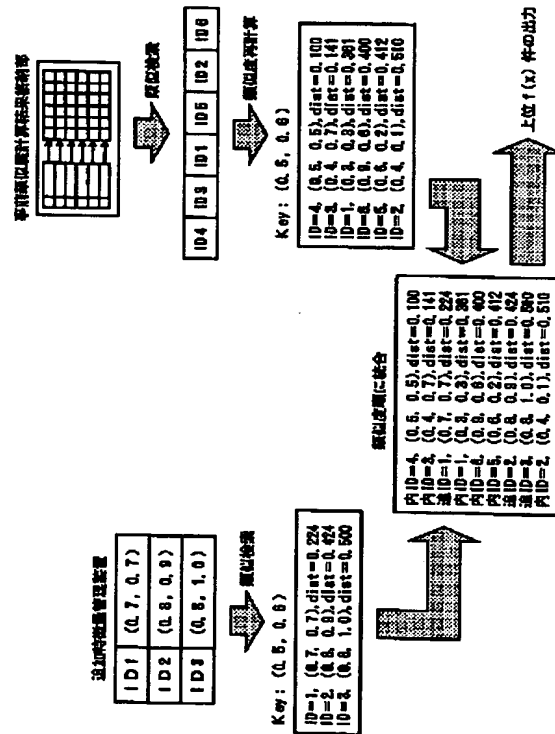
【図 3】

本発明の類似特徴量検索装置の構成図



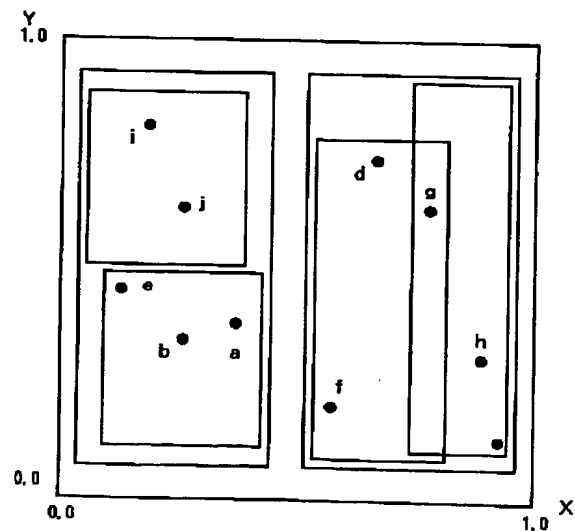
【図 9】

本発明の第 3 の実施例の追加特徴量の管理と
その類似検索の例を説明するための図



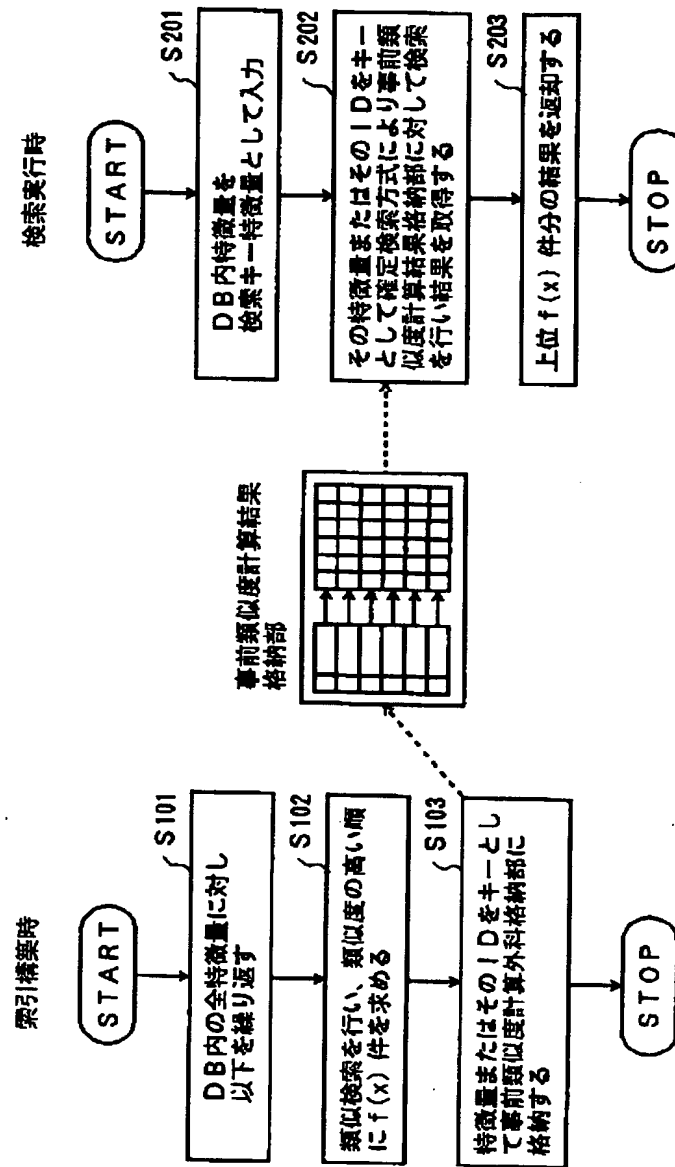
【图 1 1】

従来のR-treeの特数量空間分割を説明するための図



【図 4】

本発明の索引構築時及び検索実行時の動作を示すフローチャート

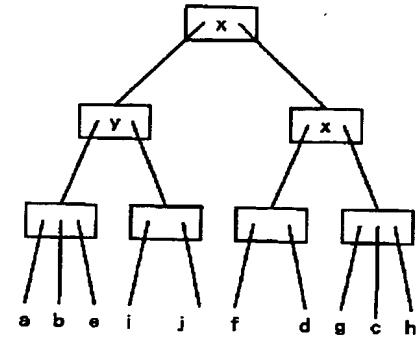
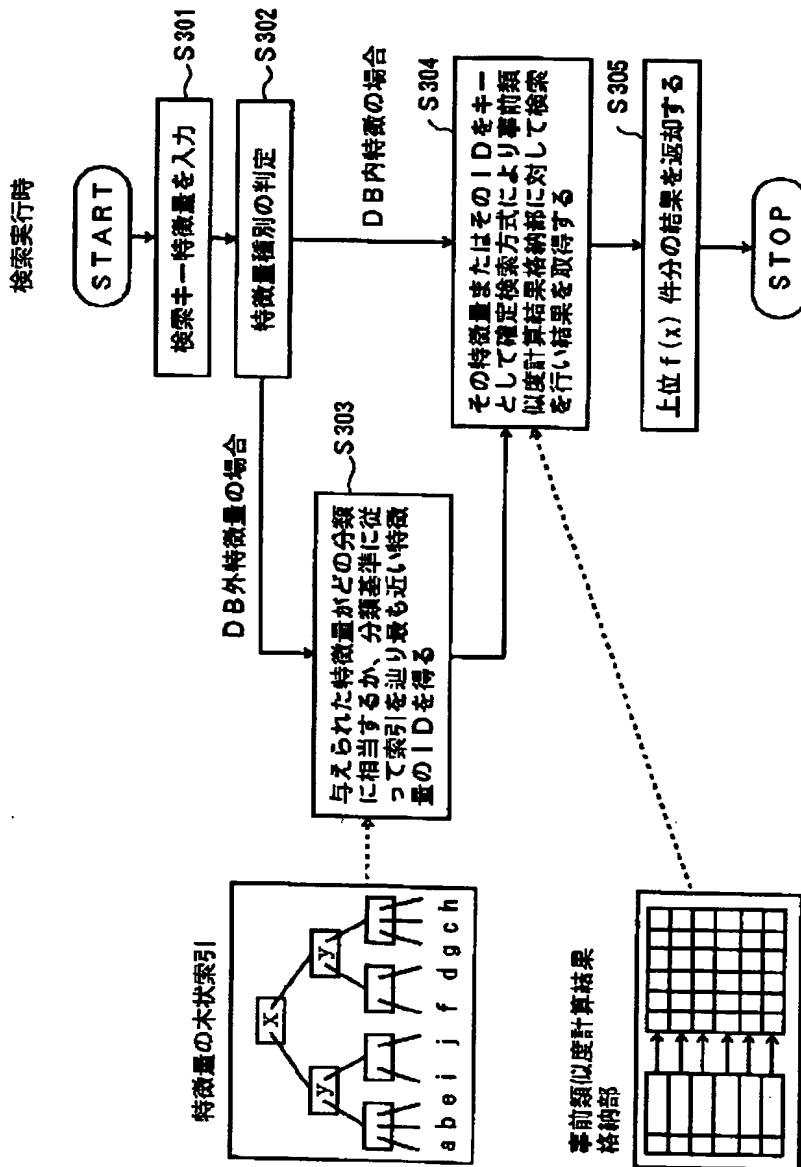


【図 5】

【図 12】

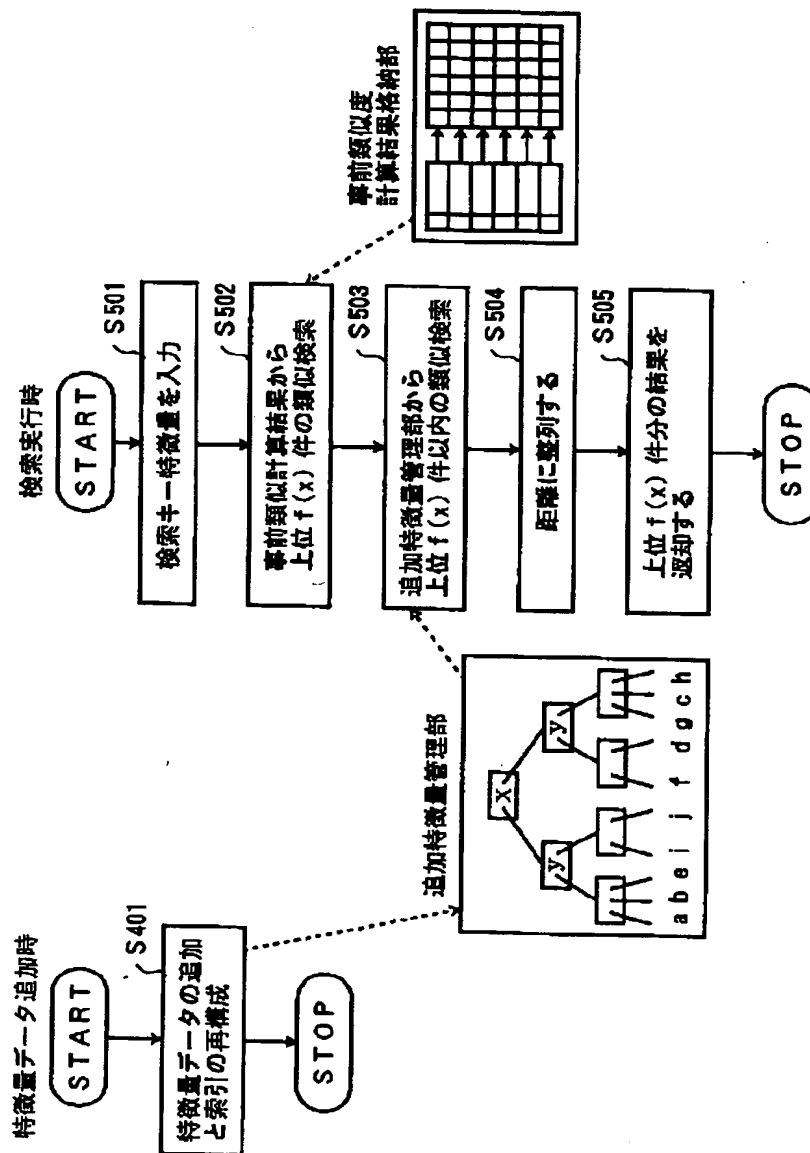
本発明の最近傍検索の処理を示すフローチャート

従来の R-tree の木構造を説明するための図



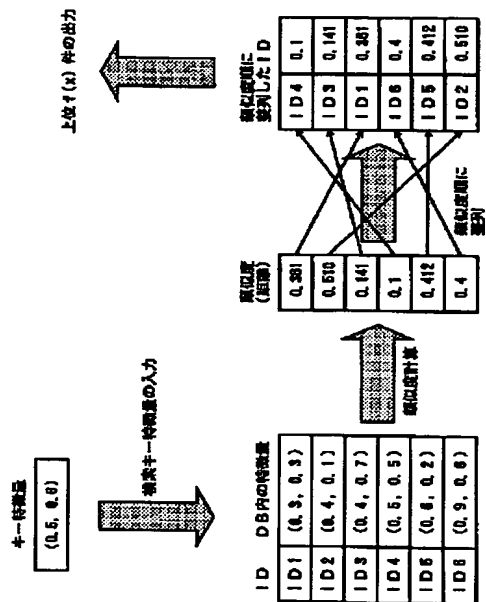
【図 6】

本発明の特徴量データの追加がある場合の
処理を示すフローチャート



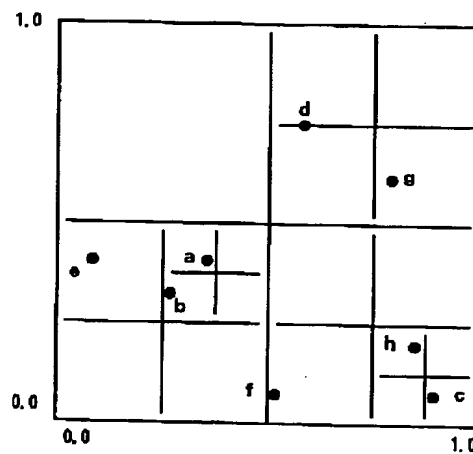
【図 10】

類似検索を説明するための図



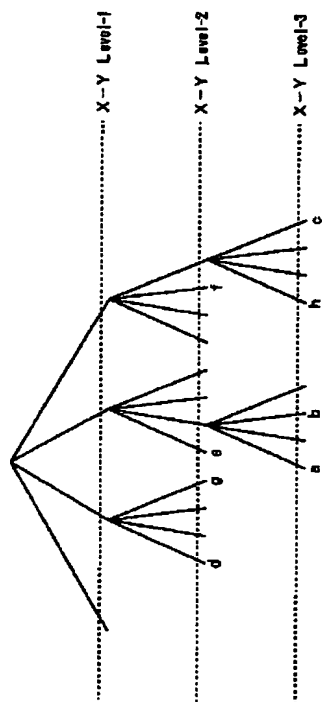
【図 13】

従来のPR-quadtreesの特徴量空間分割を説明するための図



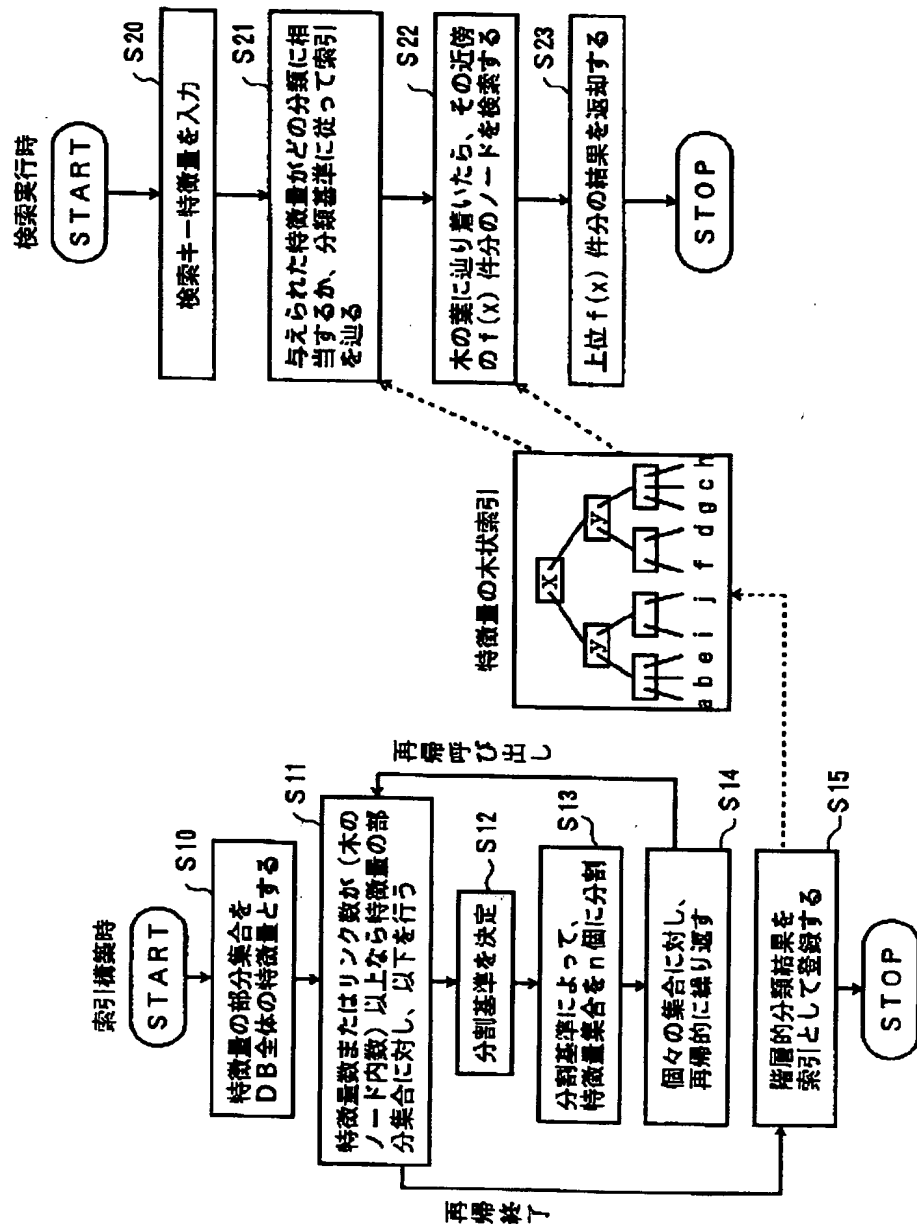
【図 14】

従来のPR-quadtreesの木構造を説明するための図



【図15】

従来の木状索引を使った類似検索のフローチャート



フロントページの続き

(72)発明者 三井 一能

東京都新宿区西新宿三丁目19番2号 日本
電信電話株式会社内

(72)発明者 串間 和彦

東京都新宿区西新宿三丁目19番2号 日本
電信電話株式会社内

F ターム(参考) 5B050 EA24 FA10 GA08
5B075 ND07 ND12 ND14 ND40 PR06
UU13 UU40
5L096 JA04 KA09

**This Page is Inserted by IFW Indexing and Scanning
Operations and is not part of the Official Record**

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☐ BLACK BORDERS
- ☐ IMAGE CUT OFF AT TOP, BOTTOM OR SIDES
- ☐ FADED TEXT OR DRAWING
- ☐ BLURRED OR ILLEGIBLE TEXT OR DRAWING
- ☐ SKEWED/SLANTED IMAGES
- ☐ COLOR OR BLACK AND WHITE PHOTOGRAPHS
- ☒ GRAY SCALE DOCUMENTS
- ☒ LINES OR MARKS ON ORIGINAL DOCUMENT
- ☐ REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY
- ☐ OTHER: _____

IMAGES ARE BEST AVAILABLE COPY.

As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.